

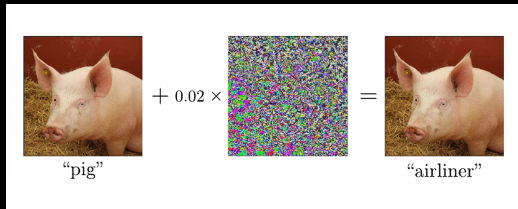
Computer Vision Layer for Robust Model

Group "Jean Ponce" :

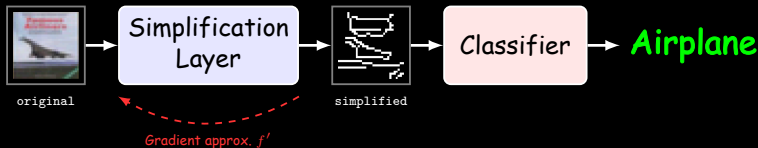
Gabriella FERNANDES MACEDO, Ruben IFRAH,
Clément ROUVROY
Dauphine-PSL, École Polytechnique, ENS-PSL

Intuition : Simplification as Defense

- ▶ **Attack** : Adds high-frequency noise (textures, details).
- ▶ **Defense** : Remove details, keep structure.



Example : Pig classified as Airliner



Canny Edge : Mechanism

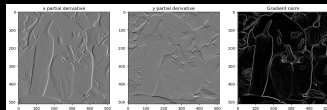
1. Gaussian Blur (Denoise)

- ▶ Removes high-freq noise.
- ▶ Prevents false detections.



2. Gradient Calculation

- ▶ Sobel operators (G_x, G_y).
- ▶ Magnitude & Direction.

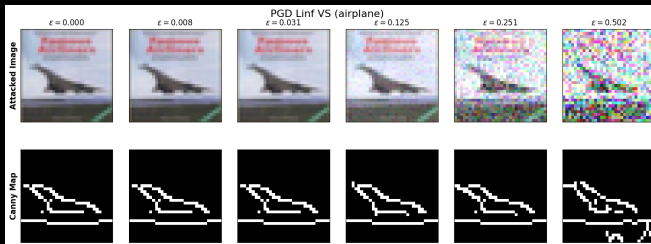


3. Thresholding : Hysteresis to connect weak edges to strong ones.

Canny Edge : Implementation

Implementation Details :

- ▶ Library : `kornia.filters.canny`
- ▶ Hyperparameter Tuning :
 - ▶ Calibrated on CIFAR-10.
 - ▶ Goal : "Is object visible?" vs "Too much noise?".



Attacked Image vs. Canny Output

Other Non-Differentiable Methods

Explored Methods :

- ▶ Quantization : Discretizing pixel values.
- ▶ Mean-Shift : Clustering-based smoothing.
- ▶ Median Filtering : Noise removal.
- ▶ Combinations : Stacking filters.

Why they failed (vs. Canny) :

- ▶ **Not destructive enough** : They preserve too much structure/linearity.
- ▶ **Ineffective against L1 attacks** : Sparse but high-magnitude noise survives these filters.
- ▶ *Canny Edge is a radical simplification that neutralizes texture attacks.*

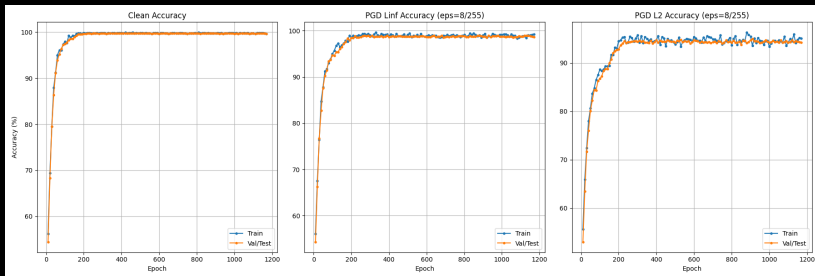
Data Augmentation

Impact on generalization :

- ▶ Without it :
Natural Acc \approx 68%.
- ▶ With it :
Natural Acc \rightarrow 100%.

Techniques :

- ▶ Random Crop : Pad 4px + Crop 32x32.
- ▶ Random Horizontal Flip : $p=0.5$.



The "Trick" : Gradient Obfuscation

To reach 197.29/200 robustness score :

Alternating Gradient in Backward Pass :

- ▶ Instead of Identity, we return $+1$ or -1 .
- ▶ Breaks the PGD optimizer (Gradient Obfuscation).

```
def backward(ctx, grad_output):  
    if call_count % 2 == 0: return 1.0  
    else: return -1.0
```

Note : Effective against PGD, but not true robustness (adaptive attacks would work).

Results

Using Canny Edge Filter + Data Augmentation :

- ▶ Nat acc : 100%
- ▶ L inf acc : 60.53% L2 acc : 89.59%

Using "the trick" :

VICTORY ROYAL

- ▶ Natural Accuracy : 100%
- ▶ Robustness Score : 197.29 / 200

Thank you for your attention

BONUS

